

# D&R: RECOVERY-BASED AI-GENERATED TEXT DETECTION VIA A SINGLE BLACK-BOX LLM CALL

000  
001  
002  
003  
004  
005 **Anonymous authors**  
006 Paper under double-blind review  
007  
008  
009  
010

## ABSTRACT

011 Large language models (LLMs) generate increasingly human-like text, raising  
012 concerns about misinformation and authenticity. Detecting AI-generated text  
013 remains challenging: existing methods often underperform, especially on short  
014 texts, require probability access unavailable in real-world black-box settings, in-  
015 incur high costs from multiple calls, or fail to generalize across models. We propose  
016 Disrupt-and-Recover (D&R), a recovery-based detection framework grounded  
017 in posterior concentration. D&R disrupts text via model-free Within-Chunk  
018 Shuffling, performs a single black-box LLM recovery, and measures seman-  
019 tic-structural recovery similarity as a proxy for concentration. This design en-  
020 sures efficiency, black-box practicality, and is theoretically supported under the  
021 concentration assumption. Extensive experiments across four datasets and six  
022 source models show that D&R achieves state-of-the-art performance, with AU-  
023 ROC 0.96 on long texts and 0.87 on short texts, surpassing the strongest base-  
024 line by +0.08 and +0.14. D&R further remains robust under source-recovery  
025 mismatch and model variation. Our code and data is available at <https://anonymous.4open.science/r/1MAdaWTy0xaod5qR>.  
026  
027

## 1 INTRODUCTION

028 Large language models (LLMs) have rapidly advanced to generate human-like text across domains  
029 such as education, news, scientific writing, and online communication. While these advances create  
030 tremendous opportunities, they also raise serious concerns about misinformation, academic integrity,  
031 and content authenticity, making reliable AI-generated text detection increasingly crucial. However,  
032 this task remains highly challenging. Real-world applications require detectors that can efficiently  
033 scale to large volumes of text with minimal overhead, for example by reducing LLM calls. They  
034 must remain robust to evolving and diverse source models while operating in black-box settings  
035 without probability access. They must also handle varied text lengths, with short texts being par-  
036 ticularly difficult. These challenges underscore the need for a detection framework that is not only  
037 accurate but also efficient, black-box practical, generalizable, and robust.  
038

039 Despite recent progress, the performance of existing AI-text detectors remains far from satisfactory,  
040 even on common long-text settings. Likelihood- and entropy-based methods (Gehrman et al., 2019;  
041 Hashimoto et al., 2019) rely on white-box access to model probabilities, making them impractical  
042 for black-box settings. Perturbation- and continuation-based methods (Bao et al., 2024; Yang et al.,  
043 2024) may improve accuracy, and rewriting-based methods (Mao et al., 2024; Park et al., 2025)  
044 avoid probability access, but all require multiple model calls, incurring high computational cost and  
045 showing instability (particularly on short texts). Supervised classifiers (OpenAI, 2019) lack gener-  
046 alization and require costly labels, while watermarking detectors (Zhao et al., 2024) heavily depend  
047 on model providers. Consequently, no existing method simultaneously delivers high performance  
048 while satisfying the demands of efficiency, black-box practicality, generalizability, and robustness.  
049

050 To address these limitations, we propose *Disrupt-and-Recover (D&R)*, a recovery-based detection  
051 framework grounded in the observation of posterior concentration: when text is disrupted in a way  
052 consistent with LLM pretraining biases, AI-generated text yields LLM-based recoveries that con-  
053 centrate more sharply around the original text  $T_{\text{orig}}$ , whereas human-written text produces more  
Shuffling (WCS), which aligns with pretraining objectives and constrains recovery to a reduced can-

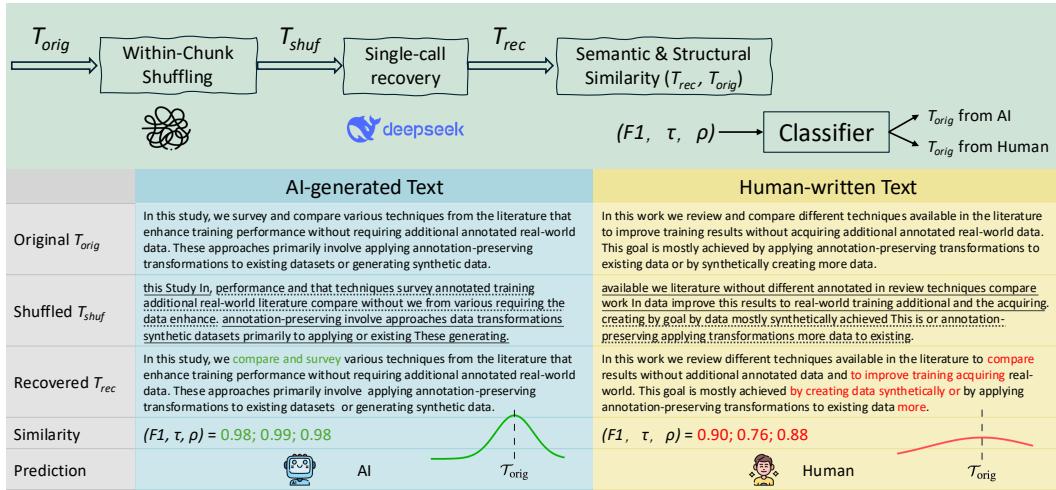


Figure 1: Illustration of D&R. **Top:** pipeline overview. **Bottom:** AI vs. Human examples. In the Shuffled  $T_{shuf}$ , individual chunks are highlighted with alternating solid and dashed underlines. AI text is recovered with fewer errors/higher similarity (green), whereas Human text shows more errors/lower similarity (red). The curves below each column schematically show recovery outcomes across examples: recoveries from AI-generated text concentrate more sharply around  $T_{orig}$ , while those from human-written text are more dispersed.

didate space; (ii) perform a single black-box recovery call, avoiding the inefficiency of multi-call methods; and (iii) compute semantic and structural recovery similarities between the recovered and original text, which serve as observable proxies for posterior concentration. These similarities are then passed to a lightweight binary classifier to output the final prediction. As illustrated in Figure 1, AI text (left) recovers with higher similarity and more concentrated outcomes near  $T_{orig}$ , while Human text (right) shows lower similarity and more dispersed recoveries. Building on this concentration assumption, we establish the theoretical foundation of D&R. We then validate its effectiveness, robustness, and usability through extensive experiments across diverse datasets, source models, and scenarios.

This paper makes the following four contributions:

- We introduce D&R, a novel detection framework that employs a model-free disruption, Within-Chunk Shuffling (WCS), conceptually aligned with the inductive biases of LLM pretraining, and ensures efficiency through a single black-box recovery call.
- We propose the Concentration Assumption and design semantic–structural recovery similarity metrics as faithful proxies for posterior concentration, providing the theoretical rationale for D&R.
- We conduct extensive experiments against eleven representative baselines across diverse datasets and source models, showing that D&R achieves state-of-the-art detection and consistently surpasses the strongest baseline (RAIDAR) by a significant margin.
- We further show that D&R remains robust and practical in challenging scenarios, including short texts, source–recovery model mismatch, and recovery–model variation.

## 2 RELATED WORK

**Zero-shot detectors.** Zero-shot detectors avoid large labeled datasets and instead exploit unsupervised signals such as probabilities or perturbations. Likelihood- and entropy-based methods (Hashimoto et al., 2019; Gehrmann et al., 2019) depend on model probability distributions, making them inherently white-box detectors. To address these limitations, recent works propose estimating black-box logits via proxy model tuning (Zeng et al., 2024), utilizing dual-model scoring (Hans et al., 2024), or analyzing intrinsic features by pre-trained models (Yu et al., 2024). Perturbation-

108 based approaches such as DetectGPT (Mitchell et al., 2023) analyze the log-likelihood curvature  
 109 of perturbed passages, with the state-of-the-art variant Fast-DetectGPT (Bao et al., 2024) improving  
 110 both accuracy and efficiency. NPR (Su et al., 2023) also leverages paraphrasing but measures  
 111 residual signals, and is less effective than Fast-DetectGPT. Continuation-based methods such as  
 112 DNA-GPT (Yang et al., 2024) truncate the text and regenerate the suffix for comparison. Rewriting-  
 113 based RAIDAR (Mao et al., 2024) generates paraphrased versions of entire passages and measures  
 114 consistency via edit distances between versions, but it requires multiple model calls, depends on  
 115 a specific paraphraser, and is vulnerable to prompt-level manipulations. Among these, RAIDAR  
 116 is most relevant to our D&R, as both rely on transformation-consistency—assessing textual con-  
 117 sistency under transformations such as paraphrasing or shuffling—recovery. In contrast to existing  
 118 generative methods (e.g., perturbation-, continuation-, and rewriting-based ones), D&R achieves  
 119 detection with only a single model call.

120 **Non-zero-shot detectors.** Non-zero-shot detectors rely on supervised discriminative models trained  
 121 with large labeled datasets. Representative examples include RoBERTa-based classifiers (Liu et al.,  
 122 2019), and the OpenAI Text Classifier (OpenAI, 2019), alongside recent frameworks utilizing multi-  
 123 level contrastive learning (Guo et al., 2024), stylistic alignment (Chen et al., 2025), or out-of-  
 124 distribution detection on human texts (Zeng et al., 2025). These approaches can achieve strong  
 125 in-domain accuracy but generalize poorly to unseen generation models and require costly labeling  
 126 and frequent retraining as LLMs evolve. Another line of work explores watermarking (Zhao et al.,  
 127 2024; Kirchenbauer et al., 2023), which embeds detectable signatures during generation. However,  
 128 watermarking depends on model providers and is unsuitable for post-hoc detection.

### 3 METHOD

132 Our D&R method follows the pipeline shown in Algorithm 1. First, we introduce a semantic-  
 133 preserving disruption, Within-Chunk Shuffling, which aligns with LLM pretraining objectives by  
 134 constraining the recovery problem to a locally permuted candidate space. Second, we perform a  
 135 single LLM call to recover the text. Next, we compute both semantic and structural similarities  
 136 between the recovered and source texts, and use these similarities as observable recoverability met-  
 137 rics. Finally, we train a binary classifier on the recoverability metrics of labeled AI-generated and  
 138 human-written texts, and apply it to obtain detection results.

---

#### Algorithm 1 D&R Pipeline

---

```

1: Input: Original text  $\mathcal{T}_{\text{orig}}$ ; black-box LLM  $\mathcal{M}$ 
2: Output: Prediction  $y \in \{\text{Human, AI}\}$ 
3: for each chunk  $c_i$  in  $\mathcal{T}_{\text{orig}}$  do
4:    $c_i^{\text{shuf}} \leftarrow \text{ShuffleTokens}(c_i)$  ▷ Apply Within-Chunk Shuffling
5: end for
6:  $\mathcal{T}_{\text{shuf}} \leftarrow \text{Join}(\{c_i^{\text{shuf}}\})$  ▷ Obtain disruption result
7:  $\mathcal{T}_{\text{rec}} \leftarrow \mathcal{M}.\text{Recover}(\mathcal{T}_{\text{shuf}})$  ▷ Single-call recovery
8:  $F_1 \leftarrow \text{SemanticSim}(\mathcal{T}_{\text{orig}}, \mathcal{T}_{\text{rec}})$  ▷ Compute semantic similarity
9:  $(\tau, \rho) \leftarrow \text{StructuralSim}(\mathcal{T}_{\text{orig}}, \mathcal{T}_{\text{rec}})$  ▷ Compute structural similarity
10:  $y \leftarrow \text{Classifier}([F_1, \tau, \rho])$  ▷ Predict label
11: return  $y$ 

```

---

153 Intuitively, recovery outcomes for AI-generated text within this constrained space tend to be highly  
 154 concentrated, whereas human-written text yields more dispersed results due to the diversity of writ-  
 155 ing processes. This concentration gap can be characterized by the notion of *posterior concentration*,  
 156 which provides the theoretical rationale behind our method and is, in practice, approximated by  
 157 recovery similarity metrics.

158 A key design in D&R is the disruption step, which determines the nature of the subsequent recov-  
 159 ery task. We adopt *Within-Chunk Shuffling* (WCS), where the original text  $\mathcal{T}_{\text{orig}}$  is segmented into  
 160 chunks by punctuation marks, and tokens within each chunk are randomly permuted while preserv-  
 161 ing chunk order. This disruption requires no model calls and can be implemented with a simple  
 random shuffling function.

162 The advantage of WCS is that it constrains recovery to a locally permuted candidate space rather  
 163 than the unconstrained generative space, closely aligning with pretraining objectives that emphasize  
 164 predicting local token orderings. As a result, recovery under WCS becomes almost effortless for the  
 165 LLM, akin to recalling the original token order, leading to recovered texts that lie very close to the  
 166 source. In distributional terms, AI-generated text tends to yield recovery outcomes that are highly  
 167 concentrated near the original text, i.e., exhibiting strong *posterior concentration*, whereas human-  
 168 written text produces more dispersed recoveries due to greater variability in writing processes.

169 Formally, as Algorithm 1 shows, we segment  $\mathcal{T}_{\text{orig}}$  into chunks by punctuation, apply a random per-  
 170 mutation to the tokens within each chunk, and then join the shuffled chunks to obtain the disrupted  
 171 text  $\mathcal{T}_{\text{shuf}}$ , which is used as the input for the subsequent recovery step.

### 173 3.1 RECOVERY WITH A SINGLE LLM CALL

175 After disruption, the shuffled text  $\mathcal{T}_{\text{shuf}}$  is passed to a large language model for recovery. We perform  
 176 this step with a *single* LLM call, where the model is prompted to restore token order and reconstruct  
 177 a coherent version of the original text  $\mathcal{T}_{\text{rec}}$ . A typical recovery prompt is:

178     The following text has its tokens shuffled within  
 179     punctuation-delimited spans. Please restore the correct  
 180     word order without adding or removing words: [INPUT].

182 This single-call design is both more efficient than multi-call approaches and well aligned with LLM  
 183 pretraining priors, as predicting local token order is a task for which pretrained models are already  
 184 highly competent. In practice, recovery can be performed either (i) via API calls to black-box  
 185 LLMs (e.g., DeepSeek-v3), or (ii) via local inference with smaller models (e.g., Mistral 7B).  
 186 Importantly, D&R achieves strong performance in both settings, demonstrating robustness and cost-  
 187 effectiveness, a property we further validate in our Recovery-Model Independence experiments (see  
 188 Section 4.3). Formally, given disrupted input  $\mathcal{T}_{\text{shuf}}$  and recovery model  $\mathcal{M}$ , the recovered text is  
 189 obtained as in Algorithm 1.

### 190 3.2 RECOVERABILITY METRICS

192 Given an original text  $\mathcal{T}_{\text{orig}}$  and its recovered counterpart  $\mathcal{T}_{\text{rec}}$ , we quantify *recoverability* using two  
 193 complementary forms of recovery similarity: semantic and structural.

194 **Semantic similarity.** We adopt *BERTScore* (Zhang et al., 2020), which measures token-level semantic  
 195 overlap by comparing contextual embeddings from a pre-trained transformer (*bert-base-uncased*,  
 196 Devlin et al., 2019). Let  $m$  and  $n$  denote the number of tokens in  $\mathcal{T}_{\text{orig}}$  and  $\mathcal{T}_{\text{rec}}$ , respectively;  $\{x_i\}_{i=1}^m$   
 197 and  $\{y_j\}_{j=1}^n$  their contextual embeddings; and  $\cos(\cdot, \cdot)$  the cosine similarity. Semantic similarity is  
 198 defined by the F1 score calculated from BScore’s precision and recall:

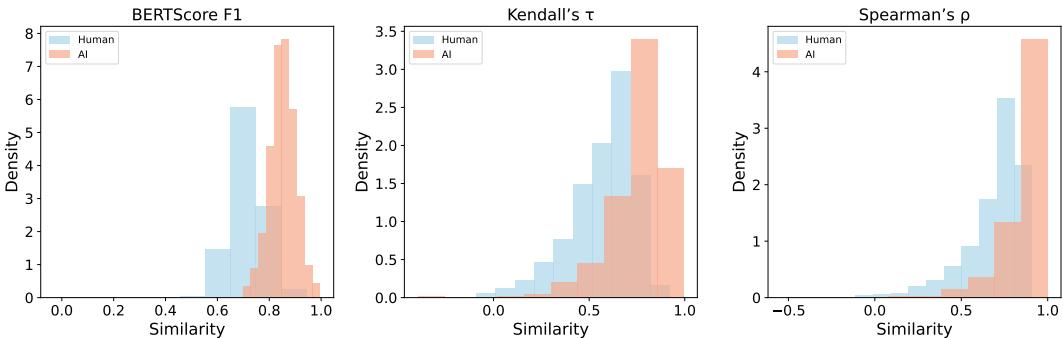
$$200 \text{Precision (P)} = \frac{1}{n} \sum_{j=1}^n \max_i \cos(x_i, y_j), \quad \text{Recall (R)} = \frac{1}{m} \sum_{i=1}^m \max_j \cos(x_i, y_j), \quad \text{F1} = \frac{2PR}{P+R}.$$

203 **Structural similarity.** We measure word-order consistency using Kendall’s  $\tau$  (Kendall, 1938; Chen  
 204 et al., 2023) and Spearman’s  $\rho$  (Spearman, 1904; Guo et al., 2025), two rank-based correlation  
 205 coefficients applied to token orderings. When  $m \neq n$  or tokens repeat, we first construct a one-  
 206 to-one alignment  $A = \{(i_k, j_k)\}_{k=1}^\ell \subseteq [m] \times [n]$  (e.g., via token-normalized Longest Common  
 207 Subsequence (LCS) with left-to-right stable matching), and then compute ranks  $r_k = i_k$  and  $s_k = j_k$   
 208 for  $k = 1, \dots, \ell$ , where  $C$  and  $D$  denote the numbers of concordant and discordant pairs among the  
 209 aligned indices.

$$210 \tau = \frac{C - D}{\frac{1}{2} \ell(\ell - 1)}, \quad \rho = 1 - \frac{6 \sum_{k=1}^\ell (r_k - s_k)^2}{\ell(\ell^2 - 1)}.$$

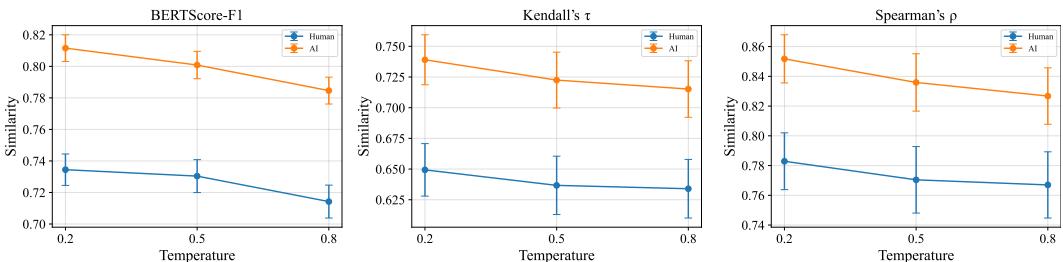
213 Semantic similarity captures fidelity of meaning, while structural similarity assesses reconstruction  
 214 of word order. Together, higher values indicate the recovered text stays closer to the source both  
 215 semantically and structurally. Thus, these complementary metrics provide observable signals of  
 recoverability.

216 **Sanity check 1.** We conducted a lightweight evaluation of our metrics on the ML-ArXiv-Papers  
 217 dataset (arXiv.org submitters, 2024) by sampling 1,000 AI-generated and 1,000 human-written texts.  
 218 For each text, we computed recovery similarity with three metrics (BERTScore F1, Kendall’s  $\tau$ ,  
 219 and Spearman’s  $\rho$ ) and plotted their distributions in Figure 2. AI-generated texts consistently  
 220 exhibit higher average similarity scores than human-written texts, revealing clear distributional gaps  
 221 across all metrics. Thus, *our recoverability metrics are effective for distinguishing AI-generated*  
 222 *from human-written text.*



234  
 235 Figure 2: Distributions of recovery similarity scores on 1,000 ML-ArXiv-Papers samples. AI-  
 236 generated texts show a clear trend toward higher scores than human-written texts across all metrics:  
 237 BERTScore F1, Kendall’s  $\tau$ , and Spearman’s  $\rho$ .

238 **Sanity check 2.** We conducted another lightweight evaluation on the ML-ArXiv-Papers  
 239 dataset (arXiv.org submitters, 2024), varying the temperature of the recovery model to control out-  
 240 put concentration, where the temperature is a decoding hyperparameter in the inference API (lower  
 241 values yield more concentrated outputs, while higher values yield more dispersed ones). As shown  
 242 in Figure 3, recovery similarity scores across all metrics decrease as temperature increases, demon-  
 243 strating a positive correlation with posterior concentration. This confirms that *our recoverability*  
 244 *metrics provide observable proxies for posterior concentration.*



254  
 255 Figure 3: Recovery similarity scores on 1,000 ML-ArXiv-Papers samples across recovery-model  
 256 temperature settings. Scores track posterior concentration: lower temperatures yield higher con-  
 257 centration and similarity, whereas higher temperatures reduce both across all metrics.

### 259 3.3 THEORETICAL ANALYSIS AND PROOF

#### 261 3.3.1 THEORETICAL RATIONALE

263 **Concentration Assumption:** Following a disruption that preserves semantics while respecting the  
 264 inductive biases of LLM pretraining (e.g., Within-Chunk Shuffling), the distribution of LLM-based  
 265 recovery outputs for AI-generated text is more concentrated in the vicinity of the original text,  
 266 whereas the distribution for human-written text is more dispersed.

267 Crucially, since posterior concentration is a distributional property that cannot be directly observed  
 268 in a single-call recovery, D&R *indirectly estimates* it from a *single recovery sample*  $T_{\text{rec}}$  by comput-  
 269 ing its similarity  $S$  to the original text  $T_{\text{orig}}$ . We prove below that **Recovery similarity is a faithful**  
**proxy for posterior concentration**, yielding a non-trivial gap between AI-and human-written texts.

270 SETUP  
271272 Let  $\mathcal{M}$  be the recovery model,  $\mathcal{T}_{\text{orig}}$  the original text,  $\mathcal{T}_{\text{shuf}}$  its WCS-disrupted version, and  $\mathcal{T}_{\text{rec}} \sim \mathcal{M}(\cdot \mid \mathcal{T}_{\text{shuf}})$  a recovered sample. Define a distance  $d(\cdot, \cdot) \geq 0$  (e.g., normalized Kendall distance) and a bounded similarity  $S(\cdot, \cdot) \in [0, 1]$  with  $S = 1$  when texts match.273  
274  
275 The posterior is  $(r, \delta)$ -concentrated if  $\Pr(d(\mathcal{T}_{\text{orig}}, \mathcal{T}_{\text{rec}}) \leq r) \geq 1 - \delta$ . Assume  $S$  is continuous in  
276  $d$  with modulus of continuity  $\omega(\cdot)$ , i.e.,  $S(\mathcal{T}_{\text{orig}}, u) \geq 1 - \omega(d(\mathcal{T}_{\text{orig}}, u))$ ,  $\forall u$ , with  $\omega(0) = 0$   
277 and  $\omega$  non-decreasing (e.g.,  $\omega(t) = Lt$ ).278  
279 **Theorem 1 (Posterior concentration  $\Rightarrow$  high recovery similarity).** If the recovery posterior is  
280  $(r, \delta)$ -concentrated, then with probability at least  $1 - \delta$ ,  $S(\mathcal{T}_{\text{orig}}, \mathcal{T}_{\text{rec}}) \geq 1 - \omega(r)$ , and consequently,

281 
$$\mathbb{E}[S(\mathcal{T}_{\text{orig}}, \mathcal{T}_{\text{rec}})] \geq (1 - \delta)(1 - \omega(r)).$$
  
282

283 *Proof.* Define  $A = \{\mathcal{T}_{\text{rec}} : d(\mathcal{T}_{\text{orig}}, \mathcal{T}_{\text{rec}}) \leq r\}$ . By posterior concentration,  $\Pr(A) \geq 1 - \delta$ . For  
284  $\mathcal{T}_{\text{rec}} \in A$ , the continuity of  $S$  gives  $S \geq 1 - \omega(r)$ . For  $\mathcal{T}_{\text{rec}} \notin A$ , we only know  $S \geq 0$ . Hence  
285  $\mathbb{E}[S] \geq (1 - \omega(r)) \Pr(A) \geq (1 - \omega(r))(1 - \delta)$ .  $\square$ 286  
287 **Theorem 2 (Non-trivial gap under Concentration Assumption).** Let  $\mathcal{T}_{\text{orig}}^{\text{AI}}$  and  $\mathcal{T}_{\text{orig}}^{\text{Human}}$  denote AI-  
288 generated and human-written texts. Suppose their recovery posteriors are  $(r_A, \delta_A)$ - and  $(r_H, \delta_H)$ -  
289 concentrated, respectively. Assume there exists  $\delta_0 > 0$  and  $\epsilon > 0$  such that the expected similarity  
290 for human text satisfies:  $\mathbb{E}[S(\mathcal{T}_{\text{orig}}^{\text{Human}}, \mathcal{T}_{\text{rec}}^{\text{Human}})] < (1 - \delta_0)(1 - \omega(r_H))$ , and  $(1 - \delta_A)(1 - \omega(r_A)) \geq$   
291  $(1 - \delta_H)(1 - \omega(r_H)) + 2\epsilon$ . Furthermore, assume the compatibility condition  $\delta_H \geq \delta_0 \geq \delta_H -$   
292  $\frac{\epsilon}{1 - \omega(r_H)}$  holds and  $\omega(r_A) \leq \omega(r_H)$ . Then

293 
$$\mathbb{E}[S(\mathcal{T}_{\text{orig}}^{\text{AI}}, \mathcal{T}_{\text{rec}}^{\text{AI}})] \geq \mathbb{E}[S(\mathcal{T}_{\text{orig}}^{\text{Human}}, \mathcal{T}_{\text{rec}}^{\text{Human}})] + \epsilon.$$
  
294

295 *Proof.* By Theorem 1,  $\mathbb{E}[S_{\text{AI}}] \geq (1 - \delta_A)(1 - \omega(r_A))$ ; by the assumption,  $\mathbb{E}[S_{\text{Human}}] < (1 - \delta_0)(1 - \omega(r_H))$ . Therefore,  $\mathbb{E}[S_{\text{AI}}] - \mathbb{E}[S_{\text{Human}}] \geq (1 - \delta_H)(1 - \omega(r_H)) + 2\epsilon - (1 - \delta_0)(1 - \omega(r_H)) = (1 - \omega(r_H))(\delta_0 - \delta_H) + 2\epsilon \geq \epsilon$ .  $\square$ 296  
297  
298 *Consequences for Metrics.* - *Kendall  $\tau$ :*  $\tau = 1 - 2d$ , so  $\omega(r) = 2r$ . - *Spearman  $\rho$ :* If at most fraction  
299  $r$  of ranks are perturbed, then  $1 - \rho \leq c_\ell r$ , hence  $\omega(r) = c_\ell r$ . - *BERTScore FI*: WCS preserves  
300 token sets; embedding drift under local permutations is bounded by  $L_{\text{sem}} r$ , hence  $\omega(r) = L_{\text{sem}} r$ .  
301302 **Takeaway.** Theorem 1 shows that posterior concentration entails high recovery similarity: as  $r \rightarrow 0$   
303 and  $\delta \rightarrow 0$ ,  $S \rightarrow 1$ . Theorem 2 shows that under the Concentration Assumption, AI texts achieve  
304 strictly higher expected recovery similarity than human texts by margin  $\epsilon$ . Thus recovery similarity  
305 is a faithful proxy for posterior concentration, providing the theoretical foundation for D&R.  
306 We provide a detailed discussion on the validity of the assumptions underlying Theorem 2 in **Appendix A.2**.  
307308 3.3.2 COMPUTATIONAL OVERHEAD.  
309310 The efficiency of D&R stems from requiring only a single black-box LLM call. Its time overhead  
311 is  $T_{\text{D\&R}} = T_{\text{shuffle}} + T_{\text{LLM}} + T_{\text{similarity}}$ , where the shuffling cost is negligible ( $T_{\text{shuffle}} \approx 0$ ) and the  
312 similarity scoring cost is much smaller than an LLM call ( $T_{\text{similarity}} \ll T_{\text{LLM}}$ ), so the overall cost is  
313 dominated by one call. In contrast, existing generative methods (e.g., perturbation-, continuation-, or  
314 rewriting-based) require multiple calls, performing  $k > 1$  queries with overhead  $T_{\text{baseline}} \approx k \cdot T_{\text{LLM}} +$   
315  $T_{\text{extra}}$ , which scales as  $O(k \cdot T_{\text{LLM}})$ . Thus, D&R lowers detection overhead to  $O(T_{\text{LLM}})$ , providing  
316 linear efficiency gains without additional assumptions. Empirical validation of these efficiency gains  
317 is provided in Appendix A.4.318 4 EXPERIMENTS  
319320 We evaluate D&R against representative baselines on long-text datasets, with ablations of recover-  
321 ability metrics, and further analyze its usability and robustness in challenging real-world scenarios  
322 including short texts, source-recovery model mismatch, and recovery LLM variation, providing a  
323 comprehensive assessment of its effectiveness.

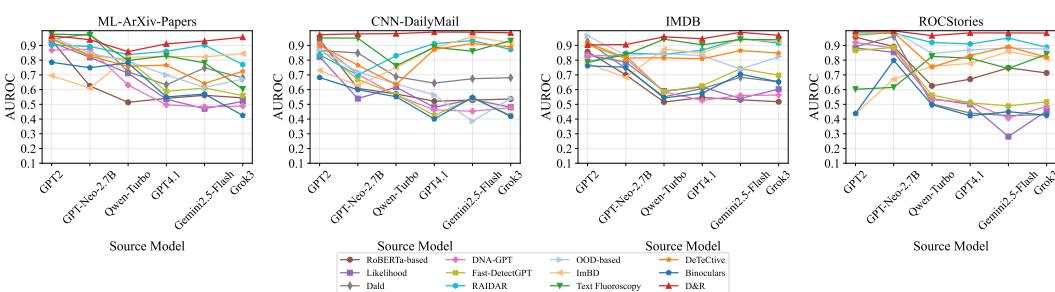
324 4.1 SETTINGS  
325

326 **Datasets and Metrics.** To evaluate the performance of the D&R method on paragraph-level AI-  
327 generated text detection, we use six publicly available datasets spanning different text lengths and  
328 domains. Based on average text length, we group them into long texts ( $>800$  words) and short texts  
329 ( $<350$  words), with dataset-wise length distributions shown in Figure 5. The long-text group in-  
330 cludes ML-ArXiv-Papers (research abstracts) (arXiv.org submitters, 2024), CNN-DailyMail (news  
331 articles) (See et al., 2017), IMDB (movie reviews) (Maas et al., 2011), and ROCStories (five-  
332 sentence stories) (Mostafazadeh et al., 2016); the short-text group includes Wikihow (instructional  
333 guides) (Sentence-Transformers, 2020), AG-News (news headlines and summaries) (Zhang et al.,  
334 2015), and Reddit (user-generated posts) (Sentence-Transformers, 2021). For each dataset, we sam-  
335 ple 1,000 human-written texts as negatives and generate AI counterparts of comparable length via  
336 paraphrasing prompts from different source models, ensuring balanced parallel datasets. To mimic  
337 diverse real-world source models, we use six widely adopted LLMs from different providers: GPT-2  
338 (Solaiman et al., 2019), GPT-Neo-2.7B (Black et al., 2021), Qwen-Turbo (Yang et al., 2025), GPT-  
339 4.1 (OpenAI, 2025), Gemini2.5-Flash (Comanici et al., 2025), and Grok3 (xAI, 2025). Detection  
340 performance is measured with the area under the ROC curve (AUROC) (Fawcett, 2006), which  
341 reflects the probability that the detector ranks an AI-generated text above a human-written one.

342 **Baselines.** We compare D&R against eleven representative baselines spanning all major families of  
343 AI-generated text detection. These include the rewriting-based RAIDAR (Mao et al., 2024), most  
344 closely related to our method; the perturbation-based Fast-DetectGPT (Bao et al., 2024), the state  
345 of the art in this family; the continuation-based DNA-GPT (Yang et al., 2024); the likelihood-based  
346 method (Hashimoto et al., 2019); recent works such as Binoculars (Hans et al., 2024), DALD (Zeng  
347 et al., 2024), and Text Fluoroscopy (Yu et al., 2024). Finally, we compare against supervised and  
348 learning-based detectors, including the RoBERTa classifier (OpenAI, 2019), Detective (Guo et al.,  
349 2024), Imitate Before Detect (Chen et al., 2025), and Human-Outlier OOD detection (Zeng et al.,  
350 2025). In the main experiments, we report results against all eleven baselines for comprehensive  
351 coverage. For additional experimental analyses, we focus on the strongest baselines and report  
352 results only against RAIDAR and Fast-DetectGPT. Baseline configurations follow their original  
353 papers, unless otherwise specified.

354 4.2 MAIN RESULTS  
355

356 **Detection Performance.** We compare D&R with eleven representative baselines on four long-text  
357 datasets across six source models, reporting AUROC performance in Table 1 and their visualization  
358 in Figure 4. For fairness, the same transformation model (DeepSeek-v3) is used for RAIDAR’s para-  
359 phrasing and D&R’s recovery. As shown in Table 1, D&R achieves the highest mean AUROC with  
360 the lowest variance ( $0.9602 \pm 0.0351$ ), substantially outperforming all baselines in both accuracy and  
361 stability. Figure 4 further illustrates that D&R’s advantage holds consistently, while baselines not  
362 only perform worse but also fluctuate with dataset and source model shifts. For instance, on ML-  
363 ArXiv-Papers, when the source model changes from GPT-2 to the more advanced Grok-3, RAIDAR  
364 drops from about 0.90 to 0.77, whereas D&R remains stable above 0.95. These results demonstrate  
365 that D&R is a robust and effective zero-shot detector for long-text scenarios.



376 Figure 4: AUROC scatter plot on four long-text datasets across six source models, complementing  
377 the averaged results in Table 1.

378 Table 1: Mean $\pm$ SD AUROC on four long-text datasets, averaged over six source models, using  
 379 DeepSeek-v3 as the recovery model. The first two entries are traditional methods, while the remain-  
 380 ing baselines represent recent state-of-the-art approaches from top venues (NeurIPS, ICML, ACL,  
 381 EMNLP, ICLR), followed by our proposed D&R. Detailed per-dataset results are provided in Ap-  
 382 pendix A.5.

| Dataset         | RoBERTa-based       | Likelihood          | Dald                | OOD-based           |
|-----------------|---------------------|---------------------|---------------------|---------------------|
| ML-ArXiv-Papers | 0.6195 $\pm$ 0.1443 | 0.6628 $\pm$ 0.1667 | 0.7838 $\pm$ 0.1441 | 0.7648 $\pm$ 0.1323 |
| CNN-DailyMail   | 0.6174 $\pm$ 0.1456 | 0.5786 $\pm$ 0.1186 | 0.7336 $\pm$ 0.0992 | 0.6203 $\pm$ 0.1703 |
| IMDB            | 0.6114 $\pm$ 0.1273 | 0.6617 $\pm$ 0.1085 | 0.6941 $\pm$ 0.1303 | 0.8392 $\pm$ 0.0773 |
| ROCStories      | 0.7675 $\pm$ 0.1182 | 0.5852 $\pm$ 0.2158 | 0.6084 $\pm$ 0.2515 | 0.9109 $\pm$ 0.0694 |
| Avg.            | 0.6539 $\pm$ 0.1498 | 0.6221 $\pm$ 0.1633 | 0.7050 $\pm$ 0.1313 | 0.7838 $\pm$ 0.1373 |
| Dataset         | ImBD                | Text Fluoroscopy    | DeTeCtive           | Binoculars          |
| ML-ArXiv-Papers | 0.7693 $\pm$ 0.0964 | 0.8266 $\pm$ 0.1367 | 0.7772 $\pm$ 0.1123 | 0.6435 $\pm$ 0.1406 |
| CNN-DailyMail   | 0.8115 $\pm$ 0.1248 | 0.8905 $\pm$ 0.0762 | 0.8295 $\pm$ 0.1116 | 0.5333 $\pm$ 0.1044 |
| IMDB            | 0.8424 $\pm$ 0.0957 | 0.8917 $\pm$ 0.0621 | 0.8452 $\pm$ 0.0403 | 0.6650 $\pm$ 0.0910 |
| ROCStories      | 0.7185 $\pm$ 0.1451 | 0.7399 $\pm$ 0.1053 | 0.8756 $\pm$ 0.0917 | 0.5054 $\pm$ 0.1487 |
| Avg.            | 0.7854 $\pm$ 0.0905 | 0.8372 $\pm$ 0.0951 | 0.8319 $\pm$ 0.0890 | 0.5868 $\pm$ 0.0962 |
| Dataset         | DNA-GPT             | Fast-DetectGPT      | RAIDAR              | D&R(ours)           |
| ML-ArXiv-Papers | 0.6400 $\pm$ 0.1708 | 0.7242 $\pm$ 0.1456 | 0.8611 $\pm$ 0.0472 | 0.9266 $\pm$ 0.0354 |
| CNN-DailyMail   | 0.5953 $\pm$ 0.1659 | 0.5838 $\pm$ 0.1556 | 0.8471 $\pm$ 0.0759 | 0.9830 $\pm$ 0.0063 |
| IMDB            | 0.6491 $\pm$ 0.1291 | 0.7277 $\pm$ 0.1075 | 0.8675 $\pm$ 0.0552 | 0.9451 $\pm$ 0.0314 |
| ROCStories      | 0.6231 $\pm$ 0.2002 | 0.6385 $\pm$ 0.1855 | 0.9323 $\pm$ 0.0482 | 0.9861 $\pm$ 0.0115 |
| Avg.            | 0.6269 $\pm$ 0.1697 | 0.6685 $\pm$ 0.1583 | 0.8770 $\pm$ 0.0657 | 0.9602 $\pm$ 0.0351 |

402 **Ablation Study.** We examine the contribution of semantic and structural recovery similarities  
 403 through an ablation study on four long-text datasets with advanced source models. As shown in  
 404 Table 7, removing semantic similarity results in the largest performance drop ( $\downarrow 28.1\%$ ), while re-  
 405 moving structural similarity also yields a substantial decrease ( $\downarrow 19.8\%$ ). The full model achieves an  
 406 AUROC of 0.9614, demonstrating that both forms of recovery similarity are indispensable and that  
 407 their combination ensures state-of-the-art accuracy and stability. We also experimentally showed  
 408 that our Within-Chunk Shuffling (WCS) is optimal compared to global or chunk-order shuffling,  
 409 effectively striking an optimal balance in the recovery task difficulty to maximize the concentration  
 410 gap, detailed analyses for these experiments are provided in Appendix A.3.2.

### 411 4.3 ANALYSIS

412 **Short-text Robustness.** Short texts are particularly challenging for AI-generated text detection,  
 413 as the limited context amplifies the distributional overlap between human and machine outputs. As  
 414 shown in Table 2, D&R achieves the highest mean AUROC with low variance ( $0.8687 \pm 0.0888$ ), sig-  
 415 nificantly outperforming RAIDAR and Fast-DetectGPT by margins of 0.14 and 0.21, respectively.  
 416 The advantage is most pronounced on earlier source models (GPT-2, GPT-Neo-2.7B), where D&R  
 417 attains near-perfect AUROC scores (around 0.99), while on stronger models performance declines  
 418 for all methods but D&R still maintains clear margins. These results underscore D&R’s consistent  
 419 superiority on short-text detection and its resilience across both weaker and stronger generators.  
 420

421 **Source-model Agnosticism (Robustness under Model Mismatch).** As transformation-  
 422 consistency based detectors, both D&R and RAIDAR rely on a transformation model to recover  
 423 or paraphrase the text generated by a source model. Although neither method requires explicit  
 424 knowledge of the source model, performance can depend on whether the source and transformation  
 425 models are the same. We therefore evaluate two cases: (i) the same-source case ( $src=trx$ ), an eas-  
 426 ier, pseudo-white-box setting in which the detector can implicitly benefit from the source model’s  
 427 distributional biases; and (ii) the different-source case ( $src \neq trx$ ), a more realistic heterogeneous  
 428 pairing. As shown in Table 3, across both settings, our D&R consistently outperforms RAIDAR.  
 429 Moreover, under the different-source condition, D&R degrades only 0.1-3.3% degradation (mean  
 430 1.9%), whereas RAIDAR drops by 4.2-14.2% (mean 9.4%). These results demonstrate that D&R  
 431 is source-agnostic: it does not rely on knowledge of the source model, remaining markedly more  
 432 robust than RAIDAR under model mismatch.

432 Table 2: AUROC performance on three Short-Text datasets across six source models. For each  
 433 dataset, results from two earlier models (GPT-2, GPT-Neo-2.7B) and four more advanced models  
 434 (Qwen-Turbo, GPT-4.1, Gemini 2.5, Grok-3) are separated by a dotted line.

| 436 | 437 | Dataset       | Source Model    | Method              |                     |                                     |
|-----|-----|---------------|-----------------|---------------------|---------------------|-------------------------------------|
|     |     |               |                 | Fast-DetectGPT      | RAIDAR              | D&R (ours)                          |
| 438 | 439 | Wikihow       | GPT2            | 0.7449              | 0.7800              | <b>0.9904</b>                       |
| 440 | 441 |               | GPT-Neo-2.7B    | 0.7936              | 0.7743              | <b>0.9987</b>                       |
| 442 | 443 |               | Qwen-Turbo      | 0.5193              | 0.5300              | <b>0.7363</b>                       |
| 444 | 445 |               | GPT4.1          | 0.4551              | 0.5700              | <b>0.7850</b>                       |
| 446 | 447 |               | Gemini2.5-Flash | 0.4370              | 0.7550              | <b>0.8517</b>                       |
| 448 | 449 |               | Grok3           | 0.4719              | 0.6950              | <b>0.7727</b>                       |
| 450 | 451 | AG-News       | GPT2            | 0.7780              | 0.7231              | <b>0.9886</b>                       |
| 452 | 453 |               | GPT-Neo-2.7B    | 0.7932              | 0.7524              | <b>0.9982</b>                       |
| 454 | 455 |               | Qwen-Turbo      | 0.7542              | 0.6850              | <b>0.7666</b>                       |
| 456 | 457 |               | GPT4.1          | 0.5819              | 0.6735              | <b>0.8202</b>                       |
| 458 | 459 |               | Gemini2.5-Flash | 0.6898              | 0.7776              | <b>0.8835</b>                       |
| 460 | 461 |               | Grok3           | 0.6439              | 0.7375              | <b>0.7963</b>                       |
| 462 | 463 | Reddit        | GPT2            | 0.7043              | 0.7649              | <b>0.9271</b>                       |
| 464 | 465 |               | GPT-Neo-2.7B    | 0.7259              | 0.7947              | <b>0.9586</b>                       |
| 466 | 467 |               | Qwen-Turbo      | 0.6852              | 0.7310              | <b>0.7502</b>                       |
| 468 | 469 |               | GPT4.1          | 0.6406              | 0.7429              | <b>0.8451</b>                       |
| 470 | 471 |               | Gemini2.5-Flash | 0.7007              | 0.7810              | <b>0.9007</b>                       |
| 472 | 473 |               | Grok3           | 0.6916              | 0.7800              | <b>0.8672</b>                       |
| 474 | 475 | Mean $\pm$ SD |                 | 0.6561 $\pm$ 0.1129 | 0.7248 $\pm$ 0.0707 | <b>0.8687<math>\pm</math>0.0888</b> |

456 Table 3: AUROC performance under same vs. different Source–Transformation Pairings. The trans-  
 457 formation model ( $trx$ ) is fixed as DeepSeek-v3. For the ‘Same’ case ( $src=trx$ ), the source model  
 458 equals the transformation model; for the ‘Different’ case ( $src\neq trx$ ), results are averaged over six  
 459 diverse source models listed in Table 1.

| 461 | 462 | Dataset         | RAIDAR             |                             | D&R (ours)          |                             |        |                    |
|-----|-----|-----------------|--------------------|-----------------------------|---------------------|-----------------------------|--------|--------------------|
|     |     |                 | Same ( $src=trx$ ) | Different ( $src\neq trx$ ) | Same ( $src=trx$ )  | Different ( $src\neq trx$ ) |        |                    |
| 463 | 464 | ML-ArXiv-Papers | 0.9475             | 0.8611                      | $\downarrow 9.1\%$  | 0.9590                      | 0.9266 | $\downarrow 3.3\%$ |
| 465 | 466 | CNN-DailyMail   | 0.9875             | 0.8471                      | $\downarrow 14.2\%$ | 0.9943                      | 0.9830 | $\downarrow 1.1\%$ |
| 467 | 468 | IMDB            | 0.9675             | 0.8675                      | $\downarrow 10.3\%$ | 0.9770                      | 0.9451 | $\downarrow 3.2\%$ |
| 469 | 470 | ROCStories      | 0.9825             | 0.9412                      | $\downarrow 4.2\%$  | 0.9869                      | 0.9865 | $\downarrow 0.1\%$ |
| 471 | 472 | Average         | 0.9712             | 0.8792                      | $\downarrow 9.4\%$  | 0.9793                      | 0.9603 | $\downarrow 1.9\%$ |

472 **Recovery-model Independence (API-based vs. Local LLMs).** We examine whether D&R de-  
 473 pends on the choice of recovery model. In addition to DeepSeek-v3 (the API-based recovery model  
 474 used in the main experiments), we evaluate Mistral-7B-Instruct-v0.3 as a locally deployed recov-  
 475 ery model. As shown in Table 4, D&R maintains strong performance (mean AUROC 0.9614 vs.  
 476 0.9359), with only  $\sim 2.5\%$  degradation when switching from a large API model to a smaller local  
 477 model. Importantly, D&R with Mistral-7B still outperforms RAIDAR even when RAIDAR relies on  
 478 the larger DeepSeek-v3 as the recovery model (data omitted for brevity). These results demon-  
 479 strate that D&R is robust across recovery-model families and scales, and remains practically deployable  
 even with smaller local models.

480 **Further Robustness and Generalization.** To thoroughly evaluate D&R’s applicability, we ex-  
 481 tended our experiments to two additional settings: (i) **Adversarial Robustness:** On the RAID  
 482 benchmark (Dugan et al., 2024), D&R retains high efficacy (AUROC 0.87) and proves resilient  
 483 against 11 varying attack categories, most notably paraphrasing. (ii) **Multilingual Generalization:**  
 484 Experiments on German, Spanish, and French confirmed that D&R generalizes effectively beyond  
 485 English, achieving  $> 0.93$  AUROC on long texts. Detailed results and analyses for these experi-  
 486 ments are provided in [Appendix A.3](#).

486  
487  
488  
Table 4: AUROC performance with two Recovery Models: DeepSeek-v3 (API-based) and Mistral-  
7B-Instruct-v0.3 (local).

| 489<br>490<br>Dataset   | Source Model    | DeepSeek-v3 (API-based)             |                     | Mistral-7B (Local)                  |
|---|-----------------|-------------------------------------|---------------------|-------------------------------------|
|   |                 | D&R                                 | RAIDAR              | D&R                                 |
| 491<br>492<br>493<br>494<br>495<br>496<br>497<br>498<br>499<br>500<br>501<br>502<br>503<br>504<br>ML-ArXiv-Papers | Qwen-Turbo      | 0.8580                              | 0.8375              | 0.8039                              |
|   | GPT4.1          | 0.9108                              | 0.8600              | 0.8656                              |
|   | Gemini2.5-Flash | 0.9299                              | 0.9025              | 0.8972                              |
|   | Grok3           | 0.9559                              | 0.7700              | 0.8157                              |
| 495<br>496<br>497<br>498<br>499<br>500<br>501<br>502<br>503<br>504<br>CNN-DailyMail                               | Qwen-Turbo      | 0.9800                              | 0.8300              | 0.9800                              |
|   | GPT4.1          | 0.9908                              | 0.9125              | 0.9844                              |
|   | Gemini2.5-Flash | 0.9901                              | 0.9325              | 0.9862                              |
|   | Grok3           | 0.9856                              | 0.8725              | 0.9784                              |
| 501<br>502<br>503<br>504<br>IMDB  | Qwen-Turbo      | 0.9584                              | 0.8400              | 0.9381                              |
|   | GPT4.1          | 0.9456                              | 0.8600              | 0.9289                              |
|   | Gemini2.5-Flash | 0.9890                              | 0.9475              | 0.9713                              |
|   | Grok3           | 0.9688                              | 0.9275              | 0.9398                              |
| 501<br>502<br>503<br>504<br>ROCStories  | Qwen-Turbo      | 0.9667                              | 0.8725              | 0.9522                              |
|   | GPT4.1          | 0.9851                              | 0.9150              | 0.9758                              |
|   | Gemini2.5-Flash | 0.9849                              | 0.9500              | 0.9818                              |
|   | Grok3           | 0.9842                              | 0.8875              | 0.9752                              |
| Mean $\pm$ SD   |                 | <b>0.9614<math>\pm</math>0.0350</b> | 0.8823 $\pm$ 0.0475 | <b>0.9359<math>\pm</math>0.0579</b> |

505  
506  
507  

## 5 CONCLUSION

508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
Disrupt-and-Recover (D&R) provides an efficient, black-box practical, and theoretically grounded  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
framework for AI-text detection, achieving state-of-the-art accuracy and robustness across diverse  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
settings, with particularly strong gains on short texts. Beyond these empirical results, D&R highlights  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568  
1569  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727  
1728  
1729  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1779  
1780  
1781  
1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1829  
1830  
1831  
1832  
1833  
1834  
1835  
1836  
1837  
1838  
1839  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889  
1889  
1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1939  
1940  
1941  
1942  
1943  
1944  
1945  
1946  
1947  
1948  
1949  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997  
1998  
1999  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2019  
2020  
2021  
2022  
2023  
2024

540 Jiaqi Chen, Xiaoye Zhu, Tianyang Liu, Ying Chen, Chen Xinhui, Yiwen Yuan, Chak Tou Leong,  
 541 Zuchao Li, Long Tang, Lei Zhang, et al. Imitate before detect: Aligning machine stylistic pref-  
 542 erence for machine-revised text detection. In *Proceedings of the AAAI Conference on Artificial*  
 543 *Intelligence*, pp. 23559–23567, 2025.

544

545 Yi Chen, Rui Wang, Haiyun Jiang, Shuming Shi, and Ruifeng Xu. Exploring the use of large  
 546 language models for reference-free text quality evaluation: An empirical study. In *Findings of*  
 547 *the Association for Computational Linguistics: IJCNLP-AACL 2023 (Findings)*, pp. 361–374.  
 548 Association for Computational Linguistics, 2023.

549

550 Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit  
 551 Dhillon, Marcel Blstein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the  
 552 frontier with advanced reasoning, multimodality, long context, and next generation agentic capa-  
 553 bilities. *arXiv preprint arXiv:2507.06261*, 2025.

554

555 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep  
 556 bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of*  
 557 *the North American chapter of the association for computational linguistics: human language*  
 558 *technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.

559

560 Liam Dugan, Alyssa Hwang, Filip Trhlík, Andrew Zhu, Josh Magnus Ludan, Hainiu Xu, Daphne  
 561 Ippolito, and Chris Callison-Burch. Raid: A shared benchmark for robust evaluation of machine-  
 562 generated text detectors. In *Proceedings of the 62nd Annual Meeting of the Association for Com-  
 563 putational Linguistics (Volume 1: Long Papers)*, pp. 12463–12492, 2024.

564

565 Tom Fawcett. An introduction to roc analysis. *Pattern Recognition Letters*, 27(8):861–874, 2006.

566

567 Sebastian Gehrmann, Hendrik Strobelt, and Alexander Rush. GLTR: Statistical detection and vi-  
 568 sualization of generated text. In *Proceedings of the 57th Annual Meeting of the Association for*  
 569 *Computational Linguistics: System Demonstrations*. Association for Computational Linguistics,  
 570 2019.

571

572 Xun Guo, Yongxin He, Shan Zhang, Ting Zhang, Wanquan Feng, Haibin Huang, and Chongyang  
 573 Ma. Detective: Detecting ai-generated text via multi-level contrastive learning. *Advances in*  
 574 *Neural Information Processing Systems*, 37:88320–88347, 2024.

575

576 Yuchen Guo, Zhicheng Dou, Huy H Nguyen, Ching-Chun Chang, Saku Sugawara, and Isao Echizen.  
 577 Measuring human involvement in ai-generated text: A case study on academic writing. *arXiv*  
 578 *preprint arXiv:2506.03501*, 2025.

579

580 Abhimanyu Hans, Avi Schwarzschild, Valeria Cherepanova, Hamid Kazemi, Aniruddha Saha,  
 581 Micah Goldblum, Jonas Geiping, and Tom Goldstein. Spotting llms with binoculars: Zero-shot  
 582 detection of machine-generated text. *arXiv preprint arXiv:2401.12070*, 2024.

583

584 Tatsunori B Hashimoto, Hugh Zhang, and Percy Liang. Unifying human and statistical evaluation  
 585 for natural language generation. *arXiv preprint arXiv:1904.02792*, 2019.

586

587 M. G. Kendall. A new measure of rank correlation. *Biometrika*, 30(1/2):81–93, 1938.

588

589 Phillip Keung, Yichao Lu, György Szarvas, and Noah A. Smith. The multilingual amazon reviews  
 590 corpus. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language*  
 591 *Processing*, 2020.

592

593 John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. A  
 594 watermark for large language models. In *Proceedings of the 40th International Conference on*  
 595 *Machine Learning*, volume 202, pp. 17061–17084. PMLR, 2023.

596

597 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike  
 598 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining  
 599 approach. *arXiv preprint arXiv:1907.11692*, 2019.

594 Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher  
 595 Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting*  
 596 *of the Association for Computational Linguistics: Human Language Technologies*, pp. 142–150.  
 597 Association for Computational Linguistics, 2011.

598 Chengzhi Mao, Carl Vondrick, Hao Wang, and Junfeng Yang. Raidar: generative AI detection via  
 599 rewriting. In *The Twelfth International Conference on Learning Representations*, 2024.

600

601 Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D Manning, and Chelsea Finn. De-  
 602 tectgpt: Zero-shot machine-generated text detection using probability curvature. In *International*  
 603 *conference on machine learning*, pp. 24950–24962. PMLR, 2023.

604 Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Van-  
 605 derwende, Pushmeet Kohli, and James Allen. A corpus and evaluation framework for deeper  
 606 understanding of commonsense stories. *arXiv preprint arXiv:1604.01696*, 2016.

607

608 OpenAI. Roberta-base-openai-detector. Hugging Face Models, 2019. URL <https://huggingface.co/roberta-base-openai-detector>.

609

610 OpenAI. Introducing gpt-4.1 in the api, 2025. URL <https://openai.com/index/gpt-4-1/>.

611

612 Hyeonchu Park, Byungjun Kim, and Bugeun Kim. Dart: An aigt detector using amr of rephrased  
 613 text. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Asso-  
 614 ciation for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*,  
 615 pp. 710–721, 2025.

616

617 Thomas Scialom, Paul-Alexis Dray, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano.  
 618 Mlsum: The multilingual summarization corpus. *arXiv preprint arXiv:2004.14900*, 2020.

619

620 Abigail See, Peter J. Liu, and Christopher D. Manning. Get to the point: Summarization with  
 621 pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for  
 622 Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics,  
 623 2017.

624

625 Sentence-Transformers. Wikihow dataset, 2020. URL <https://huggingface.co/datasets/sentence-transformers/wikihow>.

626

627 Sentence-Transformers. Reddit-title-body dataset, 2021. URL <https://huggingface.co/datasets/sentence-transformers/reddit-title-body>.

628

629 Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec  
 630 Radford, Gretchen Krueger, Jong Wook Kim, Sarah Kreps, et al. Release strategies and the social  
 631 impacts of language models. *arXiv preprint arXiv:1908.09203*, 2019.

632

633 C. Spearman. The proof and measurement of association between two things. *The American Journal  
 634 of Psychology*, 15(1):72–101, 1904.

635

636 Jinyan Su, Terry Zhuo, Di Wang, and Preslav Nakov. DetectLLM: Leveraging log rank information  
 637 for zero-shot detection of machine-generated text. In *Findings of the Association for Compu-  
 638 tational Linguistics: EMNLP 2023*. Association for Computational Linguistics, 2023.

639

640 xAI. Grok-3, 2025. URL <https://x.ai/>.

641

642 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,  
 643 Chang Gao, Chengan Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint  
 644 arXiv:2505.09388*, 2025.

645

646 Xianjun Yang, Wei Cheng, Yue Wu, Linda Ruth Petzold, William Yang Wang, and Haifeng Chen.  
 647 DNA-GPT: Divergent n-gram analysis for training-free detection of GPT-generated text. In *The  
 648 Twelfth International Conference on Learning Representations*, 2024.

649

650 Xiao Yu, Kejiang Chen, Qi Yang, Weiming Zhang, and Nenghai Yu. Text fluoroscopy: Detecting  
 651 llm-generated text through intrinsic features. In *Proceedings of the 2024 Conference on Empirical  
 652 Methods in Natural Language Processing*, pp. 15838–15846, 2024.

648 Cong Zeng, Shengkun Tang, Xianjun Yang, Yuanzhou Chen, Yiyou Sun, Zhiqiang Xu, Yao Li,  
 649 Haifeng Chen, Wei Cheng, and Dongkuan DK Xu. Dald: Improving logits-based detector without  
 650 logits from black-box llms. *Advances in Neural Information Processing Systems*, 37:54947–  
 651 54973, 2024.

652 Cong Zeng, Shengkun Tang, Yuanzhou Chen, Zhiqiang Shen, Wenchao Yu, Xujiang Zhao, Haifeng  
 653 Chen, Wei Cheng, and Zhiqiang Xu. Human texts are outliers: Detecting llm-generated texts via  
 654 out-of-distribution detection, 2025. URL <https://arxiv.org/abs/2510.08602>.

655 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evalu-  
 656 ating text generation with bert. In *International Conference on Learning Representations*, 2020.

657 Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. Character-level convolutional networks for text  
 658 classification. In *NIPS*, 2015.

659 Xuandong Zhao, Prabhanjan Vijendra Ananth, Lei Li, and Yu-Xiang Wang. Provable robust water-  
 660 marking for AI-generated text. In *The Twelfth International Conference on Learning Representa-  
 661 tions*, 2024.

## 664 A APPENDIX

### 667 A.1 THE USE OF LARGE LANGUAGE MODELS (LLMs)

668 We acknowledge the use of large language models (ChatGPT and Gemini) as assistive tools in  
 669 the preparation of this paper. Their role was strictly limited to language refinement, including  
 670 grammar correction, sentence restructuring, and style polishing. All substantive research contrib-  
 671 utions—including hypothesis formulation, experimental design and execution, result analysis, and  
 672 conclusions—are solely the work of the authors.

### 674 A.2 DISCUSSION ON THE VALIDITY OF ASSUMPTIONS IN THEOREM 2

675 In this section, we provide a detailed breakdown of the assumptions underlying Theorem 2, dis-  
 676 cussing the specific conditions under which they hold and the rationale behind them.

677 **Assumption 1 (Upper Bound of Expected Similarity for Human-Written Text).** There exist  
 678  $\delta_0 > 0$  and  $\epsilon > 0$  such that:

$$679 \mathbb{E}[S(\mathcal{T}_{orig}^{Human}, \mathcal{T}_{rec}^{Human})] < (1 - \delta_0)(1 - \omega(r_H)) \quad (1)$$

680 **When It Holds:** This assumption holds if the recovery distribution of human text satisfies two  
 681 theoretical properties:

- 682 1. *Posterior concentration:*  $\Pr(A) \geq 1 - \delta_H$ , where  $A = \{d(\mathcal{T}_{orig}^{Human}, \mathcal{T}_{rec}^{Human}) \leq r_H\}$ , as  
 683 per Theorem 1.
- 684 2. *Negligible excess similarity:* For the concentrated subset, the excess similarity is negligible,  
 685 i.e.,  $\alpha = \mathbb{E}[S | A] - (1 - \omega(r_H)) = o(1 - \omega(r_H))$ ; and for the deviated subset, the similarity  
 686 is negligible, i.e.,  $\beta = \mathbb{E}[S | A^c] = o(1 - \omega(r_H))$ .

687 **Rationale:** Human text inherently lacks the specific pretraining biases of LLMs (often exhibiting  
 688 more flexible semantics and diverse structures). Consequently, even well-recovered human text can-  
 689 not achieve the “exact consistency” typical of AI text (justifying the bound on  $\alpha$ ), while deviated  
 690 recoveries typically result in near-zero similarity (justifying  $\beta$ ). The expectation decomposition  
 691  $\mathbb{E}[S_{Human}] = \mathbb{E}[S | A] \Pr(A) + \mathbb{E}[S | A^c] \Pr(A^c)$  mathematically derives the upper bound, ensur-  
 692 ing theoretical rigor.

### 697 698 Assumption 2 (Gap in Theoretical Lower Bounds for AI-Generated Text).

$$699 (1 - \delta_A)(1 - \omega(r_A)) \geq (1 - \delta_H)(1 - \omega(r_H)) + 2\epsilon \quad (2)$$

700 **When It Holds:** This assumption holds for standard LLMs (e.g., GPT-4, Gemini) under normal  
 701 generation settings (implying a high probability of occurrence):

702 1. *Practical condition*: AI text is generated with standard parameters (e.g., temperature  $\leq 1.0$ ,  
 703 without random token insertion).

704 2. *Theoretical condition*:  $\delta_A \leq \delta_H$ ,  $r_A \leq r_H$ , and  $2\epsilon$  is less than or equal to the intrinsic  
 705 AI-human lower bound gap.

706

707 **Rationale**: LLMs are pretrained to optimize local token predictability. This objective leads to  
 708 inherently stronger posterior concentration (smaller deviation radius  $r_A$  and higher probability mass  
 709  $1 - \delta_A$ ) compared to human text, creating a natural distributional gap.

710 **Assumption 3 (Compatibility Condition)**.

712 
$$\delta_H \geq \delta_0 \geq \delta_H - \frac{\epsilon}{1 - \omega(r_H)} \quad (3)$$

713

714 **When It Holds**: This assumption holds for **all valid parameter tunings** (providing a 100% chance  
 715 of a non-empty interval):

717 1. *Practical condition*:  $\epsilon$  is set to a small conceptual margin (e.g., 0.03–0.08), aligned with  
 718 observed AI-human generation differences.

719 2. *Theoretical condition*:  $\epsilon \leq \delta_H(1 - \omega(r_H))$ . This is naturally satisfied since  $\delta_H > 0$  for  
 720 human text and  $1 - \omega(r_H) > 0$  for any reasonable radius  $r_H > 0$ .

722 **Rationale**: This interval serves to balance Assumptions 1 and 2. A value for  $\delta_0$  can always be  
 723 chosen within this range (e.g.,  $\delta_0 = \delta_H - \frac{\epsilon}{2(1 - \omega(r_H))}$ ) to strictly avoid mathematical contradictions.

725 **Summary**. The assumptions may fail only in extreme edge cases (e.g., AI text generated by non-  
 726 pretrained/random models, or human text intentionally mimicking AI patterns). However, they hold  
 727 universally in standard AI-text detection tasks. Our sanity checks confirm their practical applicabil-  
 728 ity, while the theoretical conditions ensure a high probability of occurrence in real-world scenarios.

729 **A.3 EXTENDED EXPERIMENTAL ANALYSIS**

731 In this section, we present comprehensive evaluations concerning adversarial robustness, the abla-  
 732 tion of disruption strategies, and multilingual generalization to further validate the effectiveness of  
 733 D&R.

735 **A.3.1 ROBUSTNESS AGAINST ADVERSARIAL ATTACKS**

737 To evaluate the robustness of D&R against adversarial attempts to evade detection, we utilized the  
 738 **RAID** (Dugan et al., 2024). We tested D&R against 11 diverse attack types, ranging from character-  
 739 level perturbations (e.g., homoglyphs) to high-level semantic obfuscations (e.g., paraphrasing).

740 As shown in Table 5, D&R maintains strong performance across all attack categories. Even under  
 741 **Paraphrase** attacks—typically considered the most challenging for detection—D&R maintains a  
 742 strong AUROC of 0.8210. Furthermore, for character-level attacks (e.g., Homoglyph, Zero Width  
 743 Space), performance remains robust ( $> 0.83$ ). These results indicate that our disruption-recovery  
 744 mechanism relies on intrinsic posterior concentration rather than surface-level artifacts, making it  
 745 difficult to fool via simple perturbations.

746 Table 5: Robustness of D&R on the RAID Dataset (AUROC). The method maintains high detection  
 747 performance across various attack types.

748

| Attack Type            | AUROC         | Attack Type      | AUROC  |
|------------------------|---------------|------------------|--------|
| <b>None (Clean)</b>    | <b>0.8736</b> | Homoglyph        | 0.8352 |
| Insert Paragraphs      | 0.8428        | Number           | 0.8641 |
| Alternative Spelling   | 0.8564        | Paraphrase       | 0.8210 |
| Article Deletion       | 0.8627        | Whitespace       | 0.8505 |
| Synonym                | 0.8139        | Upper/Lower      | 0.8479 |
| Perplexity Misspelling | 0.8643        | Zero Width Space | 0.8322 |

756 A.3.2 ABLATION STUDY OF WITHIN-CHUNK SHUFFLING AND RECOVERY SIMILARITY  
757758 To verify the necessity of our **Within-Chunk Shuffling (WCS)** strategy, we compared it against  
759 two alternative disruption mechanisms:760 **Global Shuffling:** Randomly shuffling all tokens in the text.  
761762 **Chunk-Order Shuffling:** Shuffling the order of chunks while keeping tokens within chunks intact.  
763764 Table 6 presents the results across four datasets.  
765766 Table 6: Ablation Study Results comparison (Avg. AUROC on 4 Datasets).  
767

| Disruption Method           | ML-ArXiv      | CNN-DM        | IMDB          | ROCStories    | Avg.          |
|-----------------------------|---------------|---------------|---------------|---------------|---------------|
| D&R (Global Shuffling)      | 0.5421        | 0.5833        | 0.5612        | 0.5390        | 0.5564        |
| D&R (Chunk-Order Shuffling) | 0.7130        | 0.7544        | 0.7205        | 0.7811        | 0.7423        |
| <b>D&amp;R (WCS - Ours)</b> | <b>0.9266</b> | <b>0.9830</b> | <b>0.9451</b> | <b>0.9861</b> | <b>0.9602</b> |

771

- 772 • **Global Shuffling:** The severe disruption destroys all semantic context, making recovery  
773 impossible for both AI and Human texts. Since both fail to be recovered, they become  
774 indistinguishable, dropping performance to random guessing ( $\sim 0.55$ ).
- 775 • **Chunk-Order Shuffling:** Preserving internal token order makes the task trivial, allowing  
776 both AI and Human texts to be recovered with high fidelity. This “ceiling effect” causes  
777 their recoverability scores to converge, significantly reducing discriminability.
- 778 • **WCS (Ours):** WCS proves to be the optimal disruption strategy. It disrupts local token  
779 order to challenge the model while preserving semantic anchors, thereby maximizing the  
780 observable “concentration gap” between AI and Human text.

782 Table 7: Ablation study of D&R by removing Semantic or Structural Recovery Similarity.  
783

| Dataset         | Source Model    | w/o SemanticSim      | w/o StructuralSim    | D&R    |
|-----------------|-----------------|----------------------|----------------------|--------|
| ML-ArXiv-Papers | Qwen-Turbo      | 0.6529               | 0.7342               | 0.8580 |
|                 | GPT4.1          | 0.6616               | 0.7268               | 0.9108 |
|                 | Gemini2.5-Flash | 0.7029               | 0.7311               | 0.9299 |
|                 | Grok3           | 0.7272               | 0.7611               | 0.9559 |
| CNN-DailyMail   | Qwen-Turbo      | 0.6567               | 0.7465               | 0.9800 |
|                 | GPT4.1          | 0.6627               | 0.7151               | 0.9908 |
|                 | Gemini2.5-Flash | 0.6955               | 0.7731               | 0.9901 |
|                 | Grok3           | 0.6653               | 0.7072               | 0.9856 |
| IMDB            | Qwen-Turbo      | 0.7112               | 0.8560               | 0.9584 |
|                 | GPT4.1          | 0.7011               | 0.8551               | 0.9456 |
|                 | Gemini2.5-Flash | 0.7013               | 0.8806               | 0.9890 |
|                 | Grok3           | 0.7159               | 0.8716               | 0.9688 |
| ROCStories      | Qwen-Turbo      | 0.6581               | 0.7158               | 0.9667 |
|                 | GPT4.1          | 0.6561               | 0.7410               | 0.9851 |
|                 | Gemini2.5-Flash | 0.8022               | 0.7910               | 0.9849 |
|                 | Grok3           | 0.6796               | 0.7256               | 0.9842 |
| Average         |                 | 0.6906 <b>↓28.1%</b> | 0.7707 <b>↓19.8%</b> | 0.9614 |

800 A.3.3 MULTILINGUAL GENERALIZATION  
801802 To demonstrate that D&R is not limited to English, we extended our experiments to **German (DE)**,  
803 **Spanish (ES)**, and **French (FR)**. We utilized the MLSUM (Scialom et al., 2020) dataset for long  
804 texts and the Amazon (Keung et al., 2020) reviews dataset for short texts, averaging results across  
805 diverse source models.  
806807 As detailed in Table 8, D&R achieves consistently high performance across all tested languages  
808 (AUROC  $> 0.93$  for long texts). Even on challenging short texts, it maintains robust performance  
809 ( $> 0.83$ ). This confirms that the principle of posterior concentration is not an artifact of English-  
810 centric training but holds across different languages.

810 Table 8: Multilingual Performance (Avg. AUROC) on long (MLSUM) and short (Amazon) texts.  
811

| 812 Language           | 813 Long Text (MLSUM) | 814 Short Text (Amazon) |
|------------------------|-----------------------|-------------------------|
| 815 German (DE)        | 0.9306                | 0.8313                  |
| 816 Spanish (ES)       | 0.9556                | 0.8604                  |
| 817 French (FR)        | 0.9377                | 0.8592                  |
| <b>818 Overall Avg</b> |                       | <b>~0.94</b>            |
| <b>819</b>             |                       | <b>~0.85</b>            |

## 820 A.4 EFFICIENCY COMPARISON

821 To quantify the practical benefits of our single-call framework, we compared the average latency  
822 and estimated cost of D&R against RAIDAR, the most competitive baseline which requires multiple  
823 generation calls.824 As shown in Table 9, D&R drastically reduces computational overhead. Specifically, the single-  
825 call design lowers the average latency from 15 seconds to 2 seconds per sample and reduces the  
826 estimated API cost from \$5 to \$0.2 per 1,000 samples. This confirms that D&R is not only accurate  
827 but also highly efficient for large-scale deployment.828 Table 9: Efficiency comparison between the multi-call baseline (RAIDAR) and our single-call  
829 method (D&R). Cost is estimated per 1,000 samples.  
830

| 831 Method                | 832 Avg Latency (s) | 833 Est. Cost (\$/1k samples) | 834 Calls per Sample |
|---------------------------|---------------------|-------------------------------|----------------------|
| 835 RAIDAR                | 836 15              | 837 \$5                       | 838 ~5 calls         |
| 839 <b>D&amp;R (Ours)</b> | 840 <b>2</b>        | 841 <b>\$0.2</b>              | 842 <b>1 call</b>    |

## 843 A.5 ADDITIONAL RESULTS

844 In this section, we present comprehensive performance data to supplement the main experimen-  
845 tal results. Tables 10, 11, 12, and 13 provide the detailed AUROC breakdown on four long-text  
846 datasets, namely ML-ArXiv-Papers, CNN-DailyMail, IMDB, and ROCStories, respectively. This  
847 table expands upon the summarized results in the main text, demonstrating D&R’s consistent su-  
848 periority across diverse source models and text domains. Table 14 reports the TPR scores at fixed FPR  
849 threshold of **1%** and **5%**.850 Table 10: AUROC on ML-ArXiv-Papers datasets across six source models.  
851

| 852 ML-ArXiv-Papers |             |                |                |                 |               |                          |
|---------------------|-------------|----------------|----------------|-----------------|---------------|--------------------------|
| 853 Source Model    | 854 RoBERTa | 855 Likelihood | 856 DNA-GPT    | 857 Fast-Detect | 858 RAIDAR    | 859 Dald                 |
| 860 GPT2            | 0.9333      | 0.9206         | 0.8679         | 0.9489          | 0.9033        | 0.9432                   |
| 861 GPT-Neo-2.7B    | 0.6245      | 0.8193         | 0.8741         | 0.8334          | 0.8933        | <b>0.9786</b>            |
| 862 Qwen-Turbo      | 0.5137      | 0.7112         | 0.6309         | 0.8038          | 0.8375        | 0.7210                   |
| 863 GPT4.1          | 0.5406      | 0.5352         | 0.4966         | 0.5878          | 0.8600        | 0.6326                   |
| 864 Gemini2.5-Flash | 0.5613      | 0.4703         | 0.4826         | 0.6115          | 0.9025        | 0.7477                   |
| 865 Grok3           | 0.5441      | 0.5206         | 0.4882         | 0.5596          | 0.7700        | 0.6798                   |
| 866 Source Model    | 867 OOD     | 868 ImBD       | 869 Binoculars | 870 Text-Flu    | 871 DeTeCtive | 872 <b>D&amp;R(ours)</b> |
| 873 GPT2            | 0.9608      | 0.6944         | 0.7856         | <b>0.9772</b>   | 0.9433        | 0.9660                   |
| 874 GPT-Neo-2.7B    | 0.8483      | 0.6103         | 0.7487         | 0.9701          | 0.8283        | 0.9390                   |
| 875 Qwen-Turbo      | 0.7935      | 0.8192         | 0.7821         | 0.7980          | 0.7581        | <b>0.8580</b>            |
| 876 GPT4.1          | 0.6995      | 0.8259         | 0.5500         | 0.8286          | 0.7666        | <b>0.9108</b>            |
| 877 Gemini2.5-Flash | 0.6156      | 0.8220         | 0.5699         | 0.7817          | 0.6416        | <b>0.9299</b>            |
| 878 Grok3           | 0.6712      | 0.8439         | 0.4248         | 0.6038          | 0.7250        | <b>0.9559</b>            |

## 879 A.6 DATASET DETAILS

880 **ML-ArXiv-Papers.** This dataset consists of abstracts from research papers in the computer science  
881 domain, particularly in machine learning, sourced from the ArXiv platform. The text is characterized

864

865

866 Table 11: AUROC on four CNN-DailyMail datasets across six source models.

| CNN-DailyMail   |         |            |            |             |           |                      |
|-----------------|---------|------------|------------|-------------|-----------|----------------------|
| Source Model    | RoBERTa | Likelihood | DNA-GPT    | Fast-Detect | RAIDAR    | Dald                 |
| GPT2            | 0.9392  | 0.8234     | 0.9126     | 0.8760      | 0.8367    | 0.8662               |
| GPT-Neo-2.7B    | 0.6084  | 0.5390     | 0.7041     | 0.8760      | 0.8367    | 0.8486               |
| Qwen-Turbo      | 0.5705  | 0.6163     | 0.5618     | 0.5675      | 0.8300    | 0.6876               |
| GPT4.1          | 0.5212  | 0.4796     | 0.4625     | 0.4240      | 0.9125    | 0.6441               |
| Gemini2.5-Flash | 0.5288  | 0.5339     | 0.4531     | 0.5419      | 0.9325    | 0.6741               |
| Grok3           | 0.5363  | 0.4799     | 0.4778     | 0.4249      | 0.8725    | 0.6808               |
| Source Model    | OOD     | ImBD       | Binoculars | Text-Flu    | DeTeCtive | <b>D&amp;R(ours)</b> |
| GPT2            | 0.9000  | 0.7303     | 0.6822     | 0.9509      | 0.9000    | <b>0.9734</b>        |
| GPT-Neo-2.7B    | 0.7666  | 0.6284     | 0.5993     | 0.9495      | 0.7666    | <b>0.9734</b>        |
| Qwen-Turbo      | 0.6346  | 0.7485     | 0.5519     | 0.7617      | 0.6346    | <b>0.9800</b>        |
| GPT4.1          | 0.8726  | 0.8852     | 0.4004     | 0.8876      | 0.8726    | <b>0.9908</b>        |
| Gemini2.5-Flash | 0.9131  | 0.9592     | 0.5471     | 0.8606      | 0.9131    | <b>0.9901</b>        |
| Grok3           | 0.8900  | 0.9175     | 0.4186     | 0.9324      | 0.8900    | <b>0.9856</b>        |

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886 Table 12: AUROC on four IMDB datasets across six source models.

| IMDB            |         |            |            |             |           |                      |
|-----------------|---------|------------|------------|-------------|-----------|----------------------|
| Source Model    | RoBERTa | Likelihood | DNA-GPT    | Fast-Detect | RAIDAR    | Dald                 |
| GPT2            | 0.8579  | 0.8340     | 0.8314     | 0.9034      | 0.7833    | <b>0.9257</b>        |
| GPT-Neo-2.7B    | 0.7000  | 0.7867     | 0.8290     | 0.8115      | 0.8467    | 0.7451               |
| Qwen-Turbo      | 0.5156  | 0.5901     | 0.5841     | 0.5862      | 0.8400    | 0.5509               |
| GPT4.1          | 0.5465  | 0.6174     | 0.5267     | 0.6257      | 0.8600    | 0.6072               |
| Gemini2.5-Flash | 0.5309  | 0.5398     | 0.5601     | 0.7422      | 0.9475    | 0.6834               |
| Grok3           | 0.5178  | 0.6027     | 0.5633     | 0.6975      | 0.9275    | 0.6525               |
| Source Model    | OOD     | ImBD       | Binoculars | Text-Flu    | DeTeCtive | <b>D&amp;R(ours)</b> |
| GPT2            | 0.9646  | 0.7780     | 0.7611     | 0.7888      | 0.9183    | 0.9036               |
| GPT-Neo-2.7B    | 0.8336  | 0.6828     | 0.7523     | 0.8357      | 0.8150    | <b>0.9056</b>        |
| Qwen-Turbo      | 0.8428  | 0.8739     | 0.5431     | 0.9424      | 0.8155    | <b>0.9584</b>        |
| GPT4.1          | 0.8319  | 0.8493     | 0.5756     | 0.9040      | 0.8091    | <b>0.9456</b>        |
| Gemini2.5-Flash | 0.7389  | 0.9442     | 0.7051     | 0.9403      | 0.8650    | <b>0.9890</b>        |
| Grok3           | 0.8234  | 0.9261     | 0.6527     | 0.9388      | 0.8483    | <b>0.9688</b>        |

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904 Table 13: AUROC on ROCStories datasets across six source models.

| ROCStories      |         |            |            |             |           |                      |
|-----------------|---------|------------|------------|-------------|-----------|----------------------|
| Source Model    | RoBERTa | Likelihood | DNA-GPT    | Fast-Detect | RAIDAR    | Dald                 |
| GPT2            | 0.8563  | 0.8833     | 0.9199     | 0.8655      | 0.9856    | 0.8915               |
| GPT-Neo-2.7B    | 0.8918  | 0.8533     | 0.8811     | 0.8862      | 0.9833    | 0.9613               |
| Qwen-Turbo      | 0.6254  | 0.5385     | 0.5312     | 0.5625      | 0.8725    | 0.5025               |
| GPT4.1          | 0.6710  | 0.4999     | 0.5132     | 0.5105      | 0.9150    | 0.4411               |
| Gemini2.5-Flash | 0.7476  | 0.2809     | 0.4066     | 0.4888      | 0.9500    | 0.4229               |
| Grok3           | 0.7131  | 0.4553     | 0.4866     | 0.5175      | 0.8875    | 0.4311               |
| Source Model    | OOD     | ImBD       | Binoculars | Text-Flu    | DeTeCtive | <b>D&amp;R(ours)</b> |
| GPT2            | 0.9956  | 0.4360     | 0.4381     | 0.6024      | 0.9691    | <b>0.9970</b>        |
| GPT-Neo-2.7B    | 0.9982  | 0.6708     | 0.7982     | 0.6152      | 0.9883    | <b>0.9988</b>        |
| Qwen-Turbo      | 0.8448  | 0.7582     | 0.4974     | 0.8257      | 0.7533    | <b>0.9667</b>        |
| GPT4.1          | 0.8689  | 0.7774     | 0.4225     | 0.8129      | 0.8316    | <b>0.9851</b>        |
| Gemini2.5-Flash | 0.8833  | 0.8575     | 0.4501     | 0.7432      | 0.8916    | <b>0.9849</b>        |
| Grok3           | 0.8745  | 0.8110     | 0.4261     | 0.8400      | 0.8199    | <b>0.9842</b>        |

917

918  
919  
920  
921 Table 14: **TPR (%) at Fixed FPR Thresholds.** Detailed performance breakdown across different  
922 source models for both long and short text settings.  
923  
924  
925  
926

| Dataset    | Metric    | Qwen | GPT  | Gemini | Grok |
|------------|-----------|------|------|--------|------|
| Long Text  | TPR@1%FPR | 74.8 | 81.8 | 90.8   | 79.5 |
|            | TPR@5%FPR | 85.4 | 89.8 | 93.8   | 94.6 |
| Short Text | TPR@1%FPR | 49.8 | 57.6 | 53.8   | 47.9 |
|            | TPR@5%FPR | 64.7 | 73.6 | 70.6   | 69.8 |

927  
928 by its professional language, rigorous structure, and strong logical coherence, representing a formal  
929 academic writing style.  
930

931 **CNN-DailyMail.** Comprising news articles from CNN and Daily Mail, this dataset is rich in factual  
932 statements and coherent narrative structures. The text is of high quality and written in accessible  
933 language, making it a common benchmark for news summarization and text generation research.

934 **IMDB.** The IMDB dataset contains a large collection of user-written movie reviews. These texts  
935 are highly subjective, feature rich linguistic expression, and convey strong sentimental polarity and  
936 personalized styles.

937 **ROCStories.** This dataset is composed of five-sentence stories centered around everyday life sce-  
938 narios. These texts exhibit clear narrative structures and causal relationships, embodying the char-  
939 acteristics of short, narrative-driven text.

940 **Wikihow.** This dataset contains texts extracted from “How-to” guides on wikiHow. The content is  
941 concise, presented in a formal style, and typically structured as clear, step-by-step instructions.

943 **AG-News.** Consisting of news headlines and short descriptions, this dataset exemplifies the style of  
944 short news text. It is highly condensed, formally structured, and logically coherent.

945 **Reddit.** This dataset is a collection of user-generated post titles and summaries from the Reddit  
946 platform. The language is colloquial and diverse in style, with a free and irregular structure that  
947 reflects the nature of social media communication.

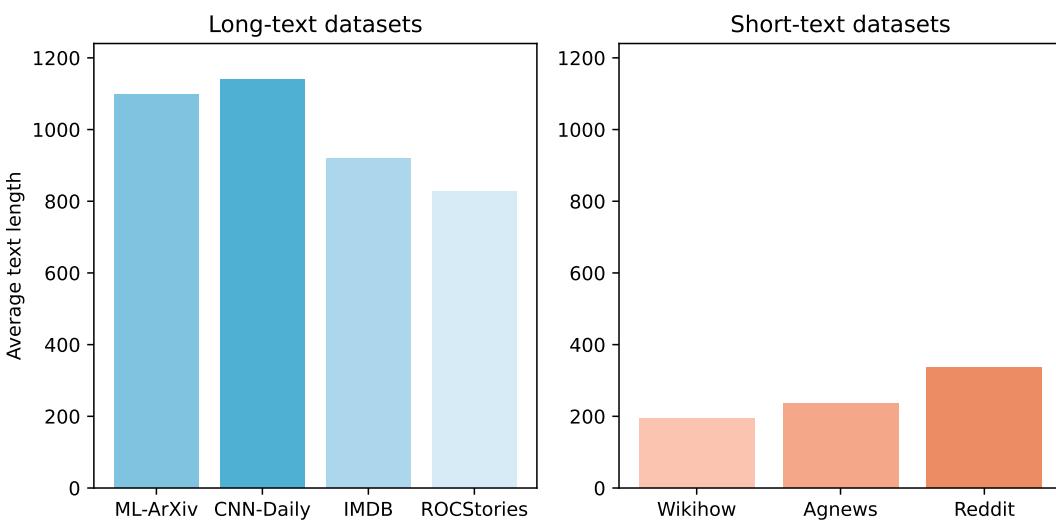


Figure 5: Average-text-length distributions for the long- and short-text datasets.